

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 05-303391

(43)Date of publication of application : 16.11.1993

(51)Int.Cl.

G10L 3/00

G10L 3/00

G10L 9/00

(21)Application number : 04-106895

(71)Applicant : SEIKO EPSON CORP

(22)Date of filing : 24.04.1992

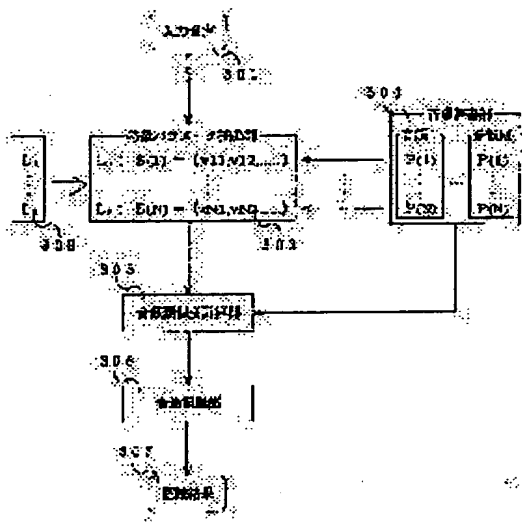
(72)Inventor : HASEGAWA HIROSHI

(54) SPEECH RECOGNITION DEVICE

(57)Abstract:

PURPOSE: To keep a high recognition rate by respectively computing feature parameter time sequences from plural frame lengths and performing respective phoneme collating.

CONSTITUTION: A feature parameter computing section 303 generates N feature parameter time sequences S1 to SN based on a frame length table 302. 16-th order LPC cepstrum coefficients are computed from 12-th order linear prediction coefficients LPC and are made to 16-dimensional feature vectors. Then, a frame is shifted to compute a feature vector. By repeating these processes, feature parameter time sequences S(j) are obtained. So far as a phoneme model is concerned, N standard patterns are made from feature parameters computed from L1 to LN frame lengths and a phoneme dictionary section 304 is constituted. A phoneme similarly computing section 305 performs a phoneme collation, outputs the phoneme with a presence probability higher than a threshold value to a phoneme recognition section 306, checks presence position, presence probability of phonemes, leaves the larger presence probability ones and outputs them.



THIS PAGE BLANK (USPTO)

2.

(19)日本国特許庁 (J P)

(12)公開特許公報 (A)

(11)特許出願公開番号

特開平5-303391

(43)公開日 平成5年(1993)11月16日

(51)Int.Cl. ⁶	識別記号	F I
G10L 3/00	515 Z 8842-5H	
	531 D 8842-5H	
	E 8842-5H	
9/00	301 A 8842-5H	

審査請求 未請求 請求項の数3 (全7頁)

(21)出願番号 特願平4-106895

(22)出願日 平成4年(1992)4月24日

(71)出願人 000002369

セイコーエプソン株式会社

東京都新宿区西新宿2丁目4番1号

(72)発明者 長谷川 浩

長野県諏訪市大和3丁目3番5号セイコー

エプソン株式会社内

(74)代理人 弁理士 鈴木 喜三郎 (外1名)

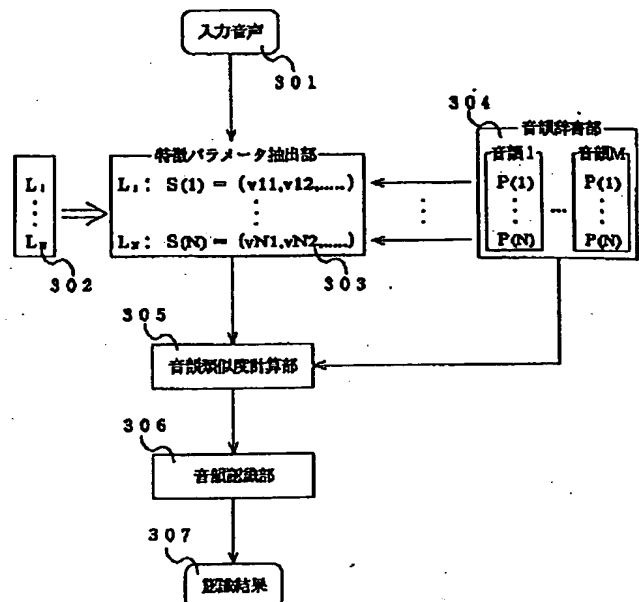
(54)【発明の名称】 音声認識装置

(57)【要約】

【目的】 入力音声から特徴パラメータを計算する際、各音韻の継続時間長の違いを反映させることにより、音韻識別率の向上をはかる。

【構成】 特徴パラメータを計算するための単位時間（フレーム）を複数個用意する、あるいは各音韻毎に用意し、各フレーム長毎に特徴パラメータ時系列を計算し、そのそれぞれに対して音韻照合を行い、最適なものを選ぶ。

【効果】 各音韻にとっての最適なフレーム長を用いることによって、音韻認識率が向上する。また、複数のフレーム長によって特徴パラメータを計算することによって、入力音声の時間的変動に対しても誤認識が少ない。



【特許請求の範囲】

【請求項1】 入力された音声データから、単位時間（1フレーム）分のデータごとに特徴パラメータを計算し、その特徴パラメータ時系列と、音韻辞書部における各音韻の標準パターンとの類似度を認識部で求め、類似度の高い音韻を認識結果とする音声認識装置において、あらかじめ特徴パラメータを計算するための単位時間長（フレーム長）を複数個（N個）備え、それぞれの音韻毎に、フレーム長 L_i により作成した標準パターンP（1）から、フレーム長 L_i により作成した標準パターンP（N）までのN個ずつの標準パターンを持つ音韻辞書部を備え、
 入力された未知音声から、前記のN個のフレーム長 $L_1 \sim L_N$ を用いてN個の特徴パラメータ時系列S（1）～S（N）を計算する特徴パラメータ計算部と、
 各音韻毎に、N個すべての特徴時系列に対して、時系列S（i）と前記標準パターンP（i）の類似度を計算し、得られたN個の類似度の最大値をその音韻の類似度として出力する類似度計算部と、
 前記類似度計算部で出力された各音韻ごとの類似度を比較し、最も類似度の高い音韻を認識結果として出力する音韻認識部と、
 を備えることを特徴とする音声認識装置。

【請求項2】 入力された音声データから、単位時間（1フレーム）分のデータごとに特徴パラメータを計算し、その特徴パラメータ時系列と、音韻辞書部における音韻の標準パターンとの類似度を音韻認識部で求め、類似度の高い音韻を認識結果とし、前記音韻認識結果から音韻系列を生成し、前記音韻系列と語彙辞書の内容の類似度を求めて単語や文節を認識する音声認識装置において、
 あらかじめ各音韻ごとに異なる長さの単位時間長（フレーム長）のデータから計算される特徴パラメータによって作成された音韻辞書部を備え、各音韻毎に最適なフレーム長およびフレームをシフトさせる時間幅を用いて特徴パラメータを計算する特徴パラメータ計算部を、各音韻毎に備えることを特徴とする音声認識装置。

【請求項3】 音韻を、その継続時間長により複数のグループに分け、個々のグループ毎に異なる長さの単位時間長（フレーム長）のデータから計算される特徴パラメータによって作成された音韻辞書部を備え、各音韻グループ毎に異なるフレーム長およびフレームをシフトさせる時間幅を用いて特徴パラメータを計算する特徴パラメータ計算部を備えることを特徴とする音声認識装置。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は音韻単位の認識に基づく音声認識装置に関するものである。

【0002】

【従来の技術】 現在考案されている音声認識装置は、そ

のほとんどが音声の特徴量の時系列に変換し、その時系列をあらかじめもっている標準パターンの時系列と比較して認識を行うというものである（図1、図2）。特徴量を計算する場合、通常数ミリ秒から数十ミリ秒を単位時間（これをフレームという）とし、1フレームの時間内では特徴量すなわち音声の波の構造は定常状態にあると近似して、LPCケプストラム等の特徴パラメータを計算する。そしてフレームをある時間だけずらして（これをフレームシフトという）、ふたたび特徴パラメータを計算する。これを繰り返すことによって特徴パラメータの時系列が得られ、これを標準パターンと比較、類似度を計算することによって認識が行われる。（たとえば特開昭61-238099）しかし音声信号は本来動的な性質をもっており、刻々とその状態は変化している。そのため1フレームの時間長（これをフレーム長という）が長すぎると、例えば/p//t//k/といった短い継続時間の子音の特徴をとらえることが難しくなる。逆にフレーム長が短いと、データ数が少なくなるため特徴パラメータの推定の精度が悪くなったり、音声にとって重要な波長の長い波の成分が見えにくくなったりする、という問題が生じてくる。

【0003】 この問題を解決するために、フレーム毎の特徴パラメータの差分ベクトルを特徴量に加え、場合によって差分ベクトルの方を重視する（特開平03-145167）等の工夫がなされているが、同じフレーム長で母音・子音の特徴パラメータをとらえているため、子音に関してはフレーム長が長すぎるため特徴量が隠され、母音に関してはフレーム長が短かすぎるためパラメータの短時間のゆらぎ等による誤認識がおこる可能性があった。これらの問題は特徴量の差分ベクトルを用いるだけでは解決できなかった。

【0004】

【発明が解決しようとする課題】 本発明の課題は、音韻認識率を大きく左右する特徴パラメータの計算方法を改善し、各音韻の認識率を向上することである。

【0005】

【課題を解決するための手段】 上記課題を解決するため、本発明の音声認識装置は、入力された音声データから、単位時間（1フレーム）分のデータごとに特徴パラメータを計算し、その特徴パラメータ時系列と、音韻辞書部における各音韻の標準パターンとの類似度を認識部で求め、類似度の高い音韻を認識結果とする音声認識装置において、あらかじめ特徴パラメータを計算するための単位時間長（フレーム長）を複数個（N個）備え、それぞれの音韻毎に、フレーム長 L_i により作成した標準パターンP（1）から、フレーム長 L_i により作成した標準パターンP（N）までのN個ずつの標準パターンを備え、入力された未知音声から、前記のN個のフレーム長 $L_1 \sim L_N$ を用いてN個の特徴パラメータ時系列S（1）～S（N）を計算する特徴パラメータ計算部を備えることを特徴とする音声認識装置。

タ計算部と、各音韻毎に、N個すべての特徴時系列に対して、時系列S(i)と前記標準パターンP(i)の類似度を計算し、得られたN個の類似度の最大値をその音韻の類似度として出力する類似度計算部と、前記類似度計算部で出力された各音韻ごとの類似度を比較し、最も類似度の高い音韻を認識結果として出力する音韻認識部と、を備えることを特徴とする。

【0006】また、入力された音声データから、単位時間(1フレーム)分のデータごとに特徴パラメータを計算し、その特徴パラメータ時系列と、音韻辞書部における音韻の標準パターンとの類似度を音韻認識部で求め、類似度の高い音韻を認識結果とし、前記音韻認識結果から音韻系列を生成し、前記音韻系列と語彙辞書の内容の類似度を求めて単語や文節を認識する音声認識装置において、あらかじめ各音韻ごとに異なる長さの単位時間長(フレーム長)のデータから計算される特徴パラメータによって作成された音韻辞書を備え、各音韻毎に最適なフレーム長およびフレームをシフトさせる時間幅を用いて特徴パラメータを計算する特徴パラメータ計算部を、各音韻毎に備えることを特徴とする。

【0007】また音韻を、その継続時間長により複数のグループに分け、個々のグループ毎に異なる長さの単位時間長(フレーム長)のデータから計算される特徴パラメータによって作成された音韻辞書を備え、各音韻グループ毎に異なるフレーム長およびフレームをシフトさせる時間幅を用いて特徴パラメータを計算する特徴パラメータ計算部を備えることを特徴とする。

【0008】

【作用】本発明は以上の構成を有するので、子音に対しては短いフレーム長で、また母音に対しては長いフレーム長で計算された特徴パラメータを用いて認識を行うことが可能となる。

【0009】

【実施例】

(実施例1) 以下本発明を実施例に基づいて詳述する。

【0010】音韻認識をおこなう場合、それぞれの音韻の平均継続時間長が問題となる。音韻の特徴は、おおきく「語頭(前の音韻の影響をうける部分)」「語中(その音韻固有の部分)」「語尾(後の音韻の影響をうける部分)」の3つに分けられる。文献(音響学会講論集1-2-14(1988-03))によると、特に短い/t//r/などの子音は語頭・語中・語尾の平均継続時間長は15ミリ秒程度しかないのに対し、母音の方はそれぞれ100ミリ秒を超える平均時間長をもつ。このように継続時間長に大きなばらつきがある様々な音素を認識するため、本発明ではあらかじめ特徴パラメータを計算するためのフレーム長を複数用意することで対処する。

【0011】図3は本発明の構成を示す図である。サンプリング周波数20kHzで16ビットで量子化された入力音声(301)から、特徴パラメータ計算部(30

3)において、特徴パラメータを計算する。特徴パラメータ計算部(303)は、あらかじめ用意されているN個のフレーム長 L_i ($i=1\dots N$)が記述されているフレーム長テーブル(302)に従って、N個の特徴パラメータ時系列 $S(1)\sim S(N)$ を生成する。この手順を説明したのが図2である。201はデジタル化された入力音声信号である。まずこの入力信号の先頭からフレーム長 L_1 (203)分のデータに注目し、この中のデータを定常状態にあるとみなして特徴パラメータを計算する。本実施例においては12次の線形予測係数LPCから16次のLPCケプストラム係数を計算して16次元の特徴ベクトル(202)とした。次にフレームをシフト(204)させ、同様に特徴ベクトルを計算する。この操作を入力信号のおわりまでくりかえすことによって、フレーム長 L_i を用いた計算した特徴パラメータ時系列 $S(i)$ が得られる。これを全てのフレーム長に関して同様に求める。その結果N個のパラメータ時系列 $S(1)$ から $S(N)$ が得られる。本実施例においては $N=5$ とし、 $L_1=3.2$, $L_2=6.4$, $L_3=12.8$, $L_4=25.6$, $L_5=51.2$ (いずれもミリ秒)とした。

【0012】一方あらかじめ個々の音韻モデルに関して、 L_1 から L_i のフレーム長から計算した特徴パラメータを用いて、N個の標準パターンを作成しておく。これは、あらかじめ発話内容と音韻の区間が既知の音声データベースを用い、それぞれのフレーム長毎に計算した特徴パラメータ時系列を、個別の隠れマルコフモデル(HMM) $P(1)\sim P(N)$ を用意してトレーニングすることによって作成した。こうして得られた音韻数 $M\times$ モデル数 N のHMMモデルにより、音韻辞書部(304)を構成した。

【0013】音韻類似度計算部(305)においては、まず各音韻毎に特徴パラメータ時系列 $S(1)$ は標準パターン $P(1)$ を用い、 $S(2)$ には $P(2)$ を用い、以下同様に $S(N)$ には $P(N)$ を用いて音韻照合を行う。そして、あらかじめ定めたしきい値を上回る存在確率をもつもののみを、音韻認識部(306)に出力する。

【0014】音韻認識部(306)では、音韻類似度計算部(305)から出力されたすべての音韻の存在位置・存在確率を調べ、存在位置が重なっているものに関しては存在確率の大きなもののみを残す。こうして得られた音韻列を認識結果(307)として出力する。

【0015】本発明により、フレーム長を固定した場合の音韻認識率と比較して、認識率の向上が認められた。

【0016】(実施例2) 図4は本発明の構成を示す図である。402は、各音韻を認識する場合最も適当なフレーム長およびフレームシフトの時間長を記述した表である。この表の値は以下のような実験によりあらかじめ求めておく。

【0017】ここでは音韻/a/に関して説明する。音声データベースから音韻/a/の発話区間を切り出し、

これを様々なフレーム長・フレームシフト時間長(N個)から計算した特徴量パラメータで、N個の標準パターンを作成する。ここではHMM音韻モデルをそれぞれの標準パターン毎に作成した。そのうち、それぞれのHMM音韻モデルを用いて、音韻/a/の認識率および/a/以外の音韻を/a/以外の音韻を/a/として認識した誤認識率を調べ、音韻/a/の識別率が最も高いものを音韻/a/の最適フレーム長・フレームシフトとした。以下全ての音韻に関して同様の実験を行い、音韻最適フレーム長テーブル(402)を作成した。

【0018】このようにして作成された最適フレーム長テーブル(402)を用いて、サンプリング周波数20kHzで16ビットで量子化された入力音声(401)から、特徴パラメータ計算部(403)において、各音韻毎に特徴パラメータを計算する。一般に最適フレーム長・シフトは音韻毎に異なるため、/a/の認識のために作成された特徴パラメータ時系列からは、音韻/a/だけが照合の対象となる。音韻認識部(405)では、音韻/a/のために計算された特徴パラメータ時系列から音韻/a/を照合し、同様に音韻/i/のための特徴パラメータ時系列から音韻/i/を照合し、以下同様に全ての音韻を、それぞれ別々に計算された特徴パラメータ時系列から照合する。こうして照合された音韻が、認識結果(406)として出力される。

【0019】本発明により、従来方法と比較して演算時間をほとんど増加させることなく、音韻認識率を向上することができた。

【0020】(実施例3)また、豊富なデータから厳密な実験を行えば、最適なフレーム長・シフトは各音韻毎に異なるが、一般には音韻を、その継続時間長に応じて3ないし4つのグループに分けて、それぞれのグループ毎に共通のフレーム長・シフトを用いれば、期待した効果が十分得られることが多い。実施例3として、音韻を「母音」「半母音・拗音」「破裂音」「その他の子音」の4グループにわけ、それぞれについてはフレーム長・シフトを共通の値とした。これにより認識率をほとんど保ったまま、演算時間を若干短縮することができた。

【0021】

【発明の効果】本発明は複数のフレーム長から特徴パラメータ時系列をそれぞれ計算し、そのそれぞれから音韻照合を行う。音声認識装置に対する入力音声は発声速度、話者の違い等により同じ音韻でもその時間的特徴は様々に変化する。従来はこういった入力に対してまず同

一フレーム長を用いて特徴量を計算し、その後の処理で入力音声のばらつきを吸収するという方法がとられている。これに対し本発明では特徴量を計算する段階で様々なフレーム長を用い、その結果最適なものを音韻ごとに選択して用いる、という手法を用いるものである。このため従来の方法と比較して入力音声の時間的特徴がばらについても、高い認識率を維持することが可能となる。

【0022】さらに、あらかじめ各音韻にとって最適なフレーム長・フレームシフト時間長を実験によって求め、その値を用いて音韻認識をすることによって、従来方法と比較してより正確な音韻認識を行うことが可能である。

【0023】また、最適フレーム長・フレームシフト時間長が類似した音韻をグループ化することによって、本手法を高速化することが可能である。

【0024】また、本発明は音声認識のデータの前処理の部分に関するものであるため、隠れマルコフモデル(HMM)等の公知のマッチング技術と組み合わせることによって、従来よりも高い認識率が得られる。

【図面の簡単な説明】

【図1】 従来の音声認識の処理を説明する図。

【図2】 特徴パラメータ抽出部の処理を説明する図。

【図3】 本発明の処理を説明する図。

【図4】 本発明の処理を説明する図。

【符号の説明】

201 入力音声

202 特徴パラメータベクトル

203 フレーム長

204 フレームシフト

301 入力音声

302 フレーム長の表

303 特徴パラメータ計算部

304 音韻辞書部

305 音韻類似度計算部

306 音韻認識部

307 認識結果

401 入力音声

402 音韻最適フレーム長テーブル

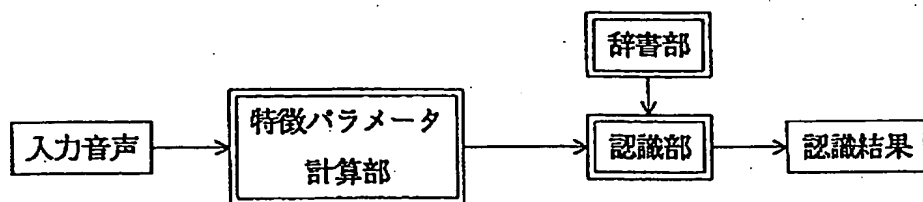
403 特徴パラメータ計算部

404 音韻辞書部

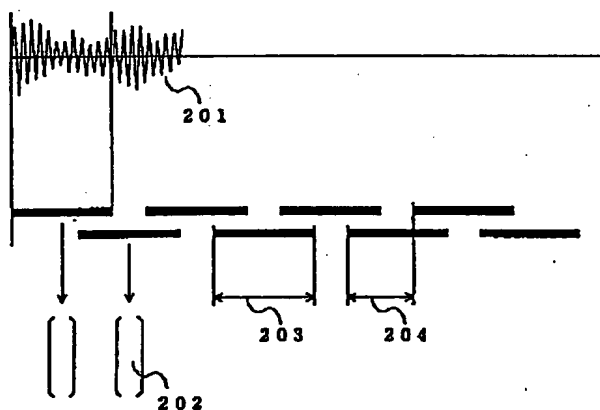
405 音韻認識部

406 認識結果

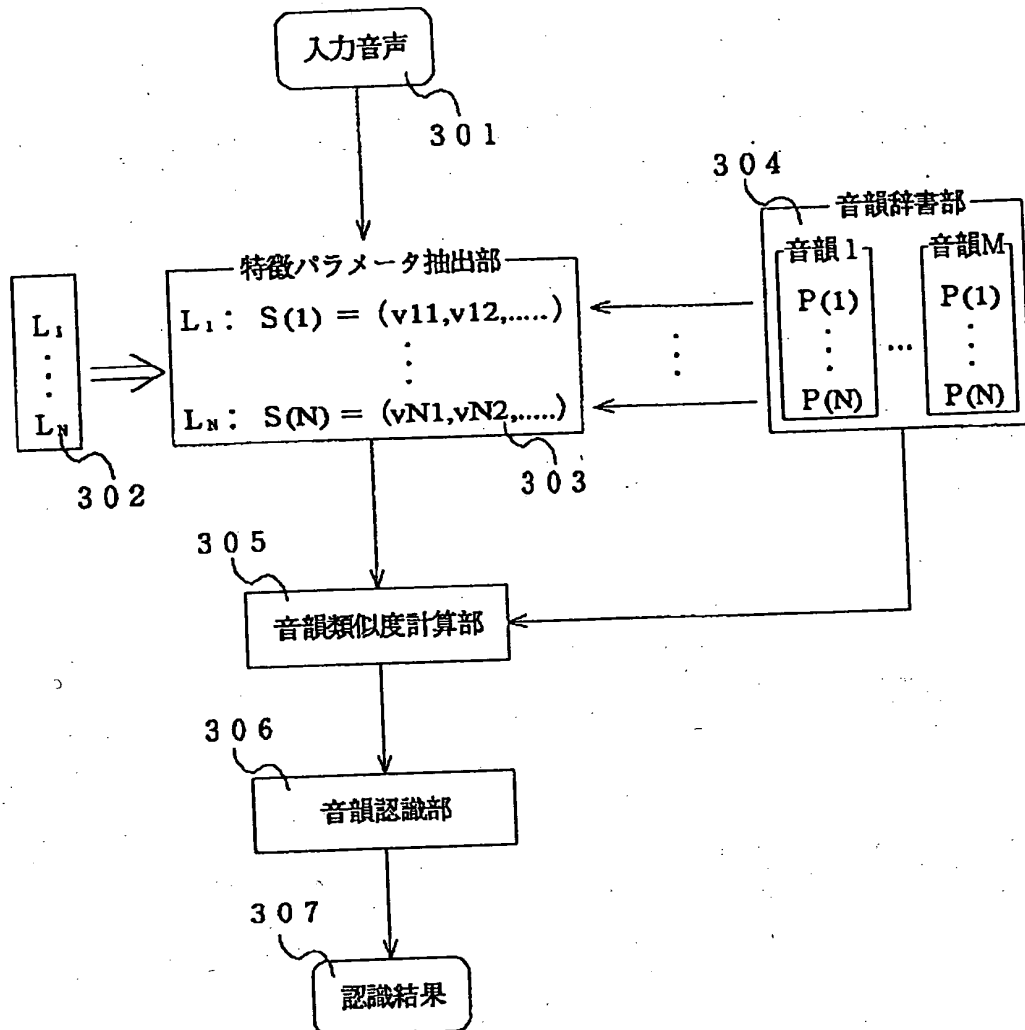
【図1】



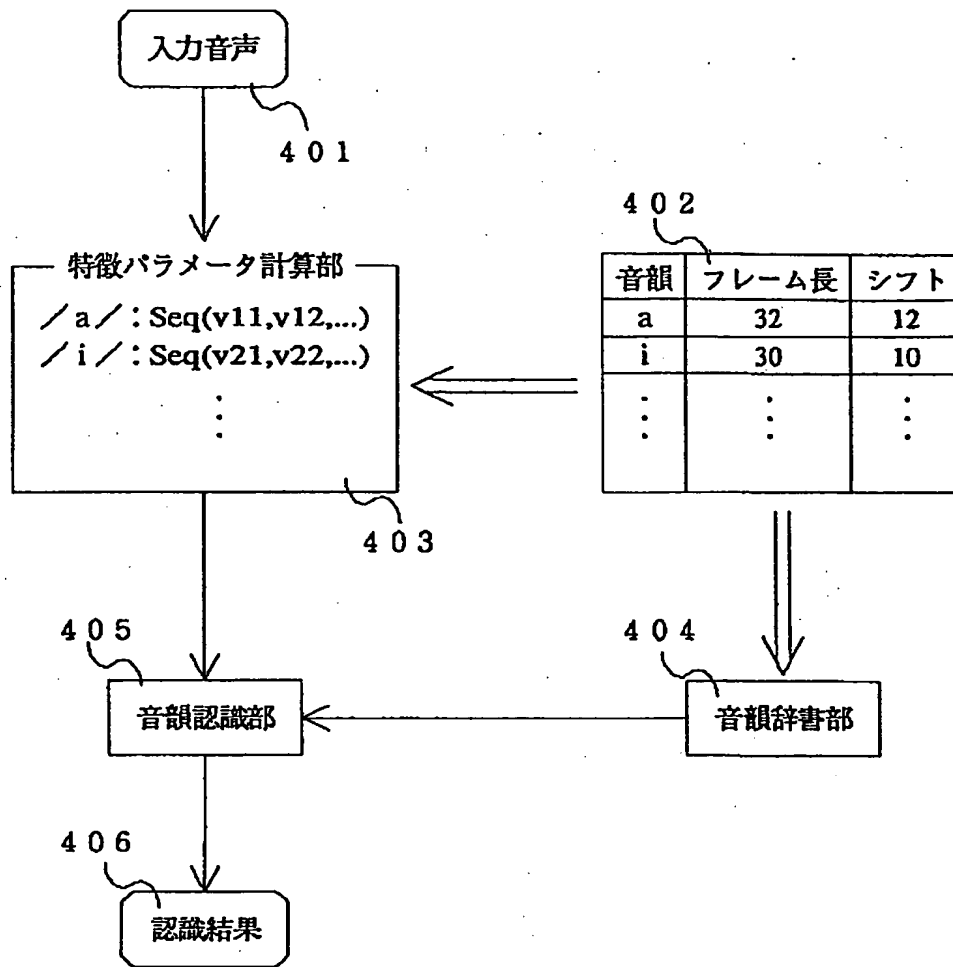
【図2】



【図3】



【図4】



THIS PAGE BLANK (USPTO)